

Automatisierung in der Kataloganreicherung

A. Wagner

`alexander.wagner@tu-ilmenau.de`

UNIVERSITÄTSBIBLIOTHEK ILMENAU

- Motivation
- Anforderungen
- Erfassung
- Texterkennung
- Upload
- Prüfung
- Zusammenfassung

Für Benutzer:

Für Benutzer:

- Suchbar
- Indexiert
- Recherchierbar

Für Benutzer:

- Suchbar
- Indexiert
- Recherchierbar

Für Bibliothek:

Für Benutzer:

- Suchbar
- Indexiert
- Recherchierbar

Für Bibliothek:

- Vereinfachung
- Beschleunigung
- Automatisierung:

Für Benutzer:

- Suchbar
- Indexiert
- Recherchierbar

Für Bibliothek:

- Vereinfachung
- Beschleunigung
- Automatisierung:
 - FTP-Transfer
 - Linkcheck
 - Transferkontrolle

1. TOC-Scan als *PDF* (bei Katalogisierung)
2. Eintrag TOC-Link im Katalog (Makro)
3. *Manueller* Upload
4. *Manueller* Linkcheck
5. *Ein* Arbeitsplatz mit Adobe Acrobat (Korrekturen)



Hohe Erkennungsrate (keine manuelle Kontrolle)

- Hohe Erkennungsrate (keine manuelle Kontrolle)
- Mehrsprachige Texte (deutsch, englisch)

- Hohe Erkennungsrate (keine manuelle Kontrolle)
- Mehrsprachige Texte (deutsch, englisch)
- Nachträgliche Texterkennung (OCR)

- Hohe Erkennungsrate (keine manuelle Kontrolle)
- Mehrsprachige Texte (deutsch, englisch)
- Nachträgliche Texterkennung (OCR)
- Flexible Integration (Geschäftsgang)

- Hohe Erkennungsrate (keine manuelle Kontrolle)
- Mehrsprachige Texte (deutsch, englisch)
- Nachträgliche Texterkennung (OCR)
- Flexible Integration (Geschäftsgang)
- *Akzeptanz* (Katalogisierung)

- Hohe Erkennungsrate (keine manuelle Kontrolle)
- Mehrsprachige Texte (deutsch, englisch)
- Nachträgliche Texterkennung (OCR)
- Flexible Integration (Geschäftsgang)
- *Akzeptanz* (Katalogisierung)
 - Einfache Erfassung
 - Geringer Zeitbedarf
 - Fehlerprüfung

- Hohe Erkennungsrate (keine manuelle Kontrolle)
- Mehrsprachige Texte (deutsch, englisch)
- Nachträgliche Texterkennung (OCR)
- Flexible Integration (Geschäftsgang)
- Akzeptanz (Katalogisierung)
 - Einfache Erfassung
 - Geringer Zeitbedarf
 - Fehlerprüfung
- Automatisierung

- Hohe Erkennungsrate (keine manuelle Kontrolle)
- Mehrsprachige Texte (deutsch, englisch)
- Nachträgliche Texterkennung (OCR)
- Flexible Integration (Geschäftsgang)
- *Akzeptanz* (Katalogisierung)
 - Einfache Erfassung
 - Geringer Zeitbedarf
 - Fehlerprüfung
- *Automatisierung*
 - Texterkennung
 - Upload
 - Linkcheck
 - Transfervergleich



Sprache: *Perl*

- Sprache: *Perl*
- *Portabel* (plattformunabhängig)

- Sprache: *Perl*
- *Portabel* (plattformunabhängig)
- Kleine Einzelscripte (flexibel)

- Sprache: *Perl*
- *Portabel* (plattformunabhängig)
- Kleine Einzelscripte (flexibel)
- Konfiguration über Variablen

- Sprache: *Perl*
- *Portabel* (plattformunabhängig)
- Kleine Einzelscripte (flexibel)
- Konfiguration über Variablen
- Arbeit zwischen zwei Verzeichnissen



Scanner:



Scanner:

- Leise (aber: Schrittmotor)
- Zügig
- Preiswert
- Einfache Software (ggf. sane)



Scanner:

- Leise (aber: Schrittmotor)
- Zügig
- Preiswert
- Einfache Software (ggf. sane)

A4-USB-Flachbettscanner



Scanner:

- Leise (aber: Schrittmotor)
- Zügig
- Preiswert
- Einfache Software (ggf. sane)

A4-USB-Flachbettscanner



Single-Page-Scan (Tiff, 400dpi, Graustufen)

- **Scanner:**
 - Leise (aber: Schrittmotor)
 - Zügig
 - Preiswert
 - Einfache Software (ggf. sane)

A4-USB-Flachbettscanner

- Single-Page-Scan (Tiff, 400dpi, Graustufen)
- Fortlaufende Numerierung (Dateiname = ppn + Nummer)

- **Scanner:**
 - Leise (aber: Schrittmotor)
 - Zügig
 - Preiswert
 - Einfache Software (ggf. sane)

A4-USB-Flachbettscanner

- Single-Page-Scan (Tiff, 400dpi, Graustufen)
- Fortlaufende Numerierung (Dateiname = ppn + Nummer)
- Einzelseiten: gemeinsamem Verzeichnis (LAN-Share)



Zeitbedarf für Scan (gegenüber PDF)

- Zeitbedarf für Scan (gegenüber PDF)
- Leiser (Schrittmotor: empfunden)

- Zeitbedarf für Scan (gegenüber PDF)
- Leiser (Schrittmotor: empfunden)
- Bessere Scanqualität

- Zeitbedarf für Scan (gegenüber PDF)
- Leiser (Schrittmotor: empfunden)
- Bessere Scanqualität
- Einfache Korrektur (kein Adobe Acrobat nötig)

- Zeitbedarf für Scan (gegenüber PDF)
- Leiser (Schrittmotor: empfunden)
- Bessere Scanqualität
- Einfache Korrektur (kein Adobe Acrobat nötig)
- OCR unabhängig vom Scanvorgang

- Zeitbedarf für Scan (gegenüber PDF)
- Leiser (Schrittmotor: empfunden)
- Bessere Scanqualität
- Einfache Korrektur (kein Adobe Acrobat nötig)
- OCR unabhängig vom Scanvorgang
- Signifikant höhere Erkennungsraten

- Zeitbedarf für Scan (gegenüber PDF)
- Leiser (Schrittmotor: empfunden)
- Bessere Scanqualität
- Einfache Korrektur (kein Adobe Acrobat nötig)
- OCR unabhängig vom Scanvorgang
- Signifikant höhere Erkennungsraten
- Nachbearbeitung bereits erfasster Dateien (transparent)

- Zeitbedarf für Scan (gegenüber PDF)
- Leiser (Schrittmotor: empfunden)
- Bessere Scanqualität
- Einfache Korrektur (kein Adobe Acrobat nötig)
- OCR unabhängig vom Scanvorgang
- Signifikant höhere Erkennungsraten
- Nachbearbeitung bereits erfasster Dateien (transparent)
- Flexibel im Geschäftsgang anzuordnen

- Viele Dateien (Jura: z. T. > 40 Dateien/Buch)
- Temporär hoher Platzbedarf (Tiff)

Vorarbeit: „Buchbinder” – createbook.pl

Vorarbeit: „Buchbinder” – createbook.pl

Abby FineReader 8 Professional

Vorarbeit: „Buchbinder” – createbook.pl

Abby FineReader 8 Professional

- Preiswert
- Hohe Erkennungsrate
- Erzeugung von Multilayer-PDF (Bild über Text)
- Unabhängig vom Eingangsformat (Nachbearbeitung alter PDFs)
- Automatische Seitenrotation (Buchscanner)

- Sonderzeichen (z. B. Mathematik)
- Speicherverwaltung FineReader (Windows-spezifisch)
- Halbautomatisch (Windows-spezifisch)

1. Eingangsverzeichnis (LAN):

<ppn>.tif, <ppn>0001.tif, <ppn>0002.tif

1. Eingangsverzeichnis (LAN):

<ppn>.tif, <ppn>0001.tif, <ppn>0002.tif

2. **Multipage-Tiff** (OCR-Input, local)

1. Eingangsverzeichnis (LAN):

<ppn>.tif, <ppn>0001.tif, <ppn>0002.tif

2. **Multipage-Tiff** (OCR-Input, local)



Ziel: ein Tiff pro ppn



Extraktion aller ppn (Dateiname)



Zusammenfügen: **tiffcp**



createbook.pl

1. Eingangsverzeichnis (LAN):

<ppn>.tif, <ppn>0001.tif, <ppn>0002.tif

2. **Multipage-Tiff** (OCR-Input, local)



Ziel: ein Tiff pro ppn



Extraktion aller ppn (Dateiname)



Zusammenfügen: **tiffcp**



createbook.pl

3. **FineReader**-Job (manuell, local):

1. Eingangsverzeichnis (LAN):

<ppn>.tif, <ppn>0001.tif, <ppn>0002.tif

2. **Multipage-Tiff** (OCR-Input, local)

● *Ziel: ein Tiff pro ppn*

● Extraktion aller ppn (Dateiname)

● Zusammenfügen: **tiffcp**

● **createbook.pl**

3. **FineReader**-Job (manuell, local):

● Lese OCR-Input

● Erzeuge Doublelayer-PDF (Upload)

1. Eingangsverzeichnis (LAN):

<ppn>.tif, <ppn>0001.tif, <ppn>0002.tif

2. **Multipage-Tiff** (OCR-Input, local)

● *Ziel: ein Tiff pro ppn*

● Extraktion aller ppn (Dateiname)

● Zusammenfügen: **tiffcp**

● **createbook.pl**

3. **FineReader**-Job (manuell, local):

● Lese OCR-Input

● Erzeuge Doublelayer-PDF (Upload)

Alles per **Batchjob** (Autostart)

Probleme

- Transferfehler zum GBV
- Fehlender Link
- Ungültiger Link (Tippfehler u. ä.)
- PPN-Änderung: manueller FTP-Zugriff

Probleme

- Transferfehler zum GBV
- Fehlender Link
- Ungültiger Link (Tippfehler u. ä.)
- PPN-Änderung: manueller FTP-Zugriff

FTP: ftp2GBV.pl

- Kompatibel
- Automatisch
- Einfach

Probleme

- Transferfehler zum GBV
- Fehlender Link
- Ungültiger Link (Tippfehler u. ä.)
- PPN-Änderung: manueller FTP-Zugriff

FTP: [ftp2GBV.pl](#)

- Kompatibel
- Automatisch
- Einfach

Linkcheck: [checkcatalog.pl](#)

- Tippfehler
- Transferfehler
- Falsche TOC

Probleme

- Transferfehler zum GBV
- Fehlender Link
- Ungültiger Link (Tippfehler u. ä.)
- PPN-Änderung: manueller FTP-Zugriff

FTP: [ftp2GBV.pl](#)

- Kompatibel
- Automatisch
- Einfach

Linkcheck: [checkcatalog.pl](#)

- Tippfehler
- Transferfehler
- Falsche TOC

Beide Schritte laufen permanent im Hintergrund



Upload/Check: Arbeitsschritte

1. Upload-Verzeichnis (local)



Upload/Check: Arbeitsschritte

1. Upload-Verzeichnis (local)
2. FTP-Transfer zum GBV

1. Upload-Verzeichnis (local)
2. FTP-Transfer zum GBV
3. Verschieben nach Check (local)

1. Upload-Verzeichnis (local)
2. FTP-Transfer zum GBV
3. Verschieben nach Check (local)
4. Abfrage des OPAC (LWP::UserAgent)



1. Upload-Verzeichnis (local)
2. FTP-Transfer zum GBV
3. Verschieben nach Check (local)
4. Abfrage des OPAC (LWP::UserAgent)
5. TOC-Link (regexp)

1. Upload-Verzeichnis (local)
2. FTP-Transfer zum GBV
3. Verschieben nach Check (local)
4. Abfrage des OPAC (LWP::UserAgent)
5. TOC-Link (regexp)
6. Download (TOC-PDF, LWP::UserAgent)

1. Upload-Verzeichnis (local)
2. FTP-Transfer zum GBV
3. Verschieben nach Check (local)
4. Abfrage des OPAC (LWP::UserAgent)
5. TOC-Link (regexp)
6. Download (TOC-PDF, LWP::UserAgent)
7. Bitvergleich (Download und lokale Kopie, File::Compare)

1. Upload-Verzeichnis (local)
2. FTP-Transfer zum GBV
3. Verschieben nach Check (local)
4. Abfrage des OPAC (LWP::UserAgent)
5. TOC-Link (regexp)
6. Download (TOC-PDF, LWP::UserAgent)
7. Bitvergleich (Download und lokale Kopie, File::Compare)
 - OK: verschiebe von Check nach Backup (LAN)
 - Fehler: verschiebe von Check nach Upload (local)
 - Falsche/Fehlende Links: Error-Log

Archivierung:

-  Tiff-Input (multipage)
-  Korrekt übertragene PDFs

Archivierung:

- Tiff-Input (multipage)
- Korrekt übertragene PDFs

Archivhandling:

- Archivierung: 2 Wochen
- [cleanup.pl](#)

Komplettes Setup seit ca. 1 Jahr im Einsatz

Komplettes Setup seit ca. 1 Jahr im Einsatz

- ✓ Beschleunigung des Arbeitsprozesses
- ✓ Gute Akzeptanz bei den Mitarbeiterinnen
- ✓ Gute Auslastung des OCR-PC
- ✓ Niedriger Wartungsaufwand
- ✓ Komponentennutzung auch in anderen Bereichen
- ✓ **PicaPM**: gute Perl-SRU-Schnittstelle für CBS/LBS